

# Variational Methods for Latent Variable Problems (part 2)

---

Ryan Giordano (for Johns Hopkins Biostats BLAST working group)

Nov 2021

Massachusetts Institute of Technology

Outline for today:

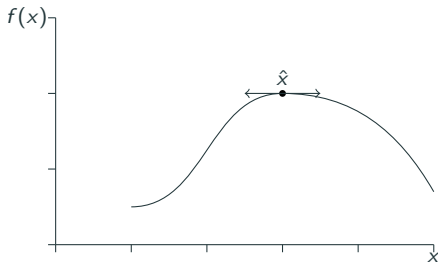
- What counts as variational inference?
- Kullback-Leibler (KL) divergence and “standard” variational inference
- The classical EM algorithm as a special case of variational inference
- Variational inference as a generalization of the EM algorithm
- Some examples of VI in practice

# What counts as variational inference?

Lots of very different procedures go by the name “variational inference.” I propose an (idiosyncratic) encompassing definition based on the use cases and the name:

**Variational inference is inference using optimization.**

Think “calculus of variations:” an optimum  $\hat{x} = \underset{\theta}{\operatorname{argmax}} f(x)$  is characterized by  $df/dx|_{\hat{x}} = 0$ , i.e. where small variations in  $\hat{x}$  result in no changes to the value of  $f(\hat{x})$ .

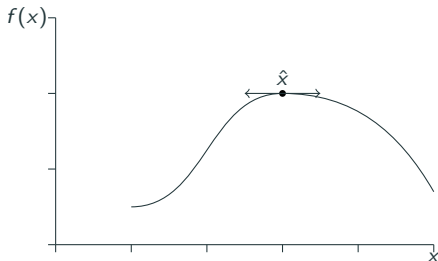


# What counts as variational inference?

Lots of very different procedures go by the name “variational inference.” I propose an (idiosyncratic) encompassing definition based on the use cases and the name:

**Variational inference is inference using optimization.**

Think “calculus of variations:” an optimum  $\hat{x} = \underset{\theta}{\operatorname{argmax}} f(x)$  is characterized by  $df/dx|_{\hat{x}} = 0$ , i.e. where small variations in  $\hat{x}$  result in no changes to the value of  $f(\hat{x})$ .



**Exercise:** By this definition, which of these are VI?

- The maximum likelihood estimator (MLE).
- The Laplace approximation to a Bayesian posterior.
- Markov chain Monte Carlo (MCMC).

# What counts as variational inference?

A more common definition of VI is the following.

Suppose we have a random variable  $\xi$  and a distribution  $p(\xi)$  that we want to know.

Let  $y$  denote data and  $\theta$  a parameter. Examples:

- The variable is  $\theta$ , and we wish to know the posterior  $p(\theta|y)$  (Bayes)
- The variable is  $y$ , and we wish to know  $p(y)$  (MLE)
- The variable is  $y$ , and we wish to know the map  $\theta \mapsto p(y|\theta) = \int p(y, z|\theta) dz$  (marginal MLE)

# What counts as variational inference?

A more common definition of VI is the following.

Suppose we have a random variable  $\xi$  and a distribution  $p(\xi)$  that we want to know.

Let  $y$  denote data and  $\theta$  a parameter. Examples:

- The variable is  $\theta$ , and we wish to know the posterior  $p(\theta|y)$  (Bayes)
- The variable is  $y$ , and we wish to know  $p(y)$  (MLE)
- The variable is  $y$ , and we wish to know the map  $\theta \mapsto p(y|\theta) = \int p(y, z|\theta) dz$  (marginal MLE)

Let  $\mathcal{Q}$  be some class of distributions which may or may not contain  $p(\xi)$ .

**Variational inference finds the distribution in  $\mathcal{Q}$  closest to  $p$  according to some measure of “divergence” between distributions:**

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} D(q, p).$$

# What counts as variational inference?

A more common definition of VI is the following.

Suppose we have a random variable  $\xi$  and a distribution  $p(\xi)$  that we want to know.

Let  $y$  denote data and  $\theta$  a parameter. Examples:

- The variable is  $\theta$ , and we wish to know the posterior  $p(\theta|y)$  (Bayes)
- The variable is  $y$ , and we wish to know  $p(y)$  (MLE)
- The variable is  $y$ , and we wish to know the map  $\theta \mapsto p(y|\theta) = \int p(y, z|\theta) dz$  (marginal MLE)

Let  $\mathcal{Q}$  be some class of distributions which may or may not contain  $p(\xi)$ .

**Variational inference finds the distribution in  $\mathcal{Q}$  closest to  $p$  according to some measure of “divergence” between distributions:**

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} D(q, p).$$

The most common choice of “divergence” is the **Kullback-Leibler** (KL) divergence, though other choices are possible (e.g. Li and Turner [2016], Liu and Wang [2016], Ambrogioni et al. [2018]).

# KL divergence

The KL divergence is defined as:

$$\text{KL}(q||p) := \mathbb{E}_{q(\xi)} [\log q(\xi)] - \mathbb{E}_{q(\xi)} [\log p(\xi)]$$

Some key attributes of KL divergence:

- $\text{KL}(q||p) \geq 0$
- $\text{KL}(q||p) = 0 \Rightarrow p = q$
- $\text{KL}(q||p) \neq \text{KL}(p||q)$
- $\text{KL}(q||p)$  is a “strict” measure of closeness [Gibbs and Su, 2002]



# KL divergence

The KL divergence is defined as:

$$\text{KL}(q||p) := \mathbb{E}_{q(\xi)} [\log q(\xi)] - \mathbb{E}_{q(\xi)} [\log p(\xi)]$$

Some key attributes of KL divergence:

- $\text{KL}(q||p) \geq 0$
- $\text{KL}(q||p) = 0 \Rightarrow p = q$
- $\text{KL}(q||p) \neq \text{KL}(p||q)$
- $\text{KL}(q||p)$  is a “strict” measure of closeness [Gibbs and Su, 2002]

Why use KL divergence?

**Phony answer:** The KL divergence has an information theoretic interpretation [Kullback and Leibler, 1951].

**Real answer:** Mathematical convenience (normalizing constants pop out).

# KL divergence

The KL divergence is defined as:

$$\text{KL}(q||p) := \mathbb{E}_{q(\xi)} [\log q(\xi)] - \mathbb{E}_{q(\xi)} [\log p(\xi)]$$

Some key attributes of KL divergence:

- $\text{KL}(q||p) \geq 0$
- $\text{KL}(q||p) = 0 \Rightarrow p = q$
- $\text{KL}(q||p) \neq \text{KL}(p||q)$
- $\text{KL}(q||p)$  is a “strict” measure of closeness [Gibbs and Su, 2002]

Why use KL divergence?

**Phony answer:** The KL divergence has an information theoretic interpretation [Kullback and Leibler, 1951].

**Real answer:** Mathematical convenience (normalizing constants pop out).

**Example: The MLE minimizes KL divergence.** Suppose that  $x_n \stackrel{iid}{\sim} p(\cdot)$ , and  $q(\cdot|\theta) \in \mathcal{Q}$  is a (possibly misspecified) parameteric family of data distributions. Then

$$\begin{aligned} \operatorname{argmin}_{\theta} \text{KL}(p||q) &= \operatorname{argmin}_{\theta} \left( -\mathbb{E}_{p(x_1)} [\log q(x_1|\theta)] + \mathbb{E}_{p(x_1)} [\log p(x_1)] \right) \\ &= \operatorname{argmax}_{\theta} \mathbb{E}_{p(x_1)} [\log q(x_1|\theta)] \approx \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{n=1}^N \log q(x_n|\theta) = \hat{\theta} \text{ (the MLE)}. \end{aligned}$$

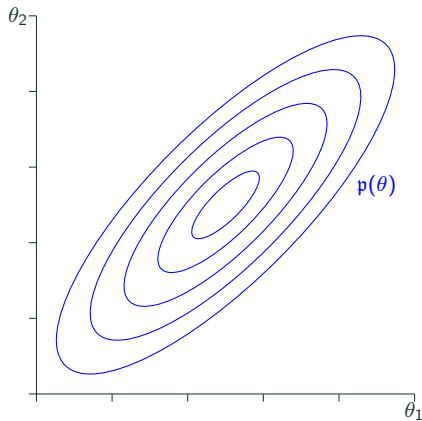
# KL divergence exercises

$$\text{KL}(q(\theta) || p(\theta)) = - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{All bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



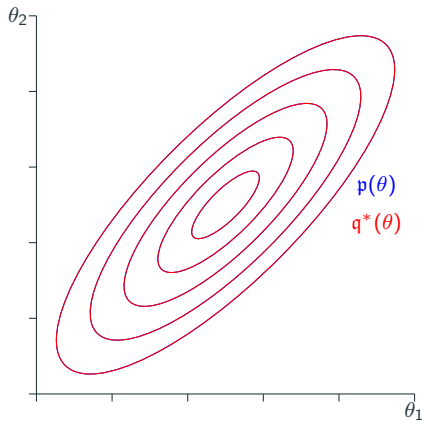
# KL divergence exercises

$$\text{KL}(q(\theta) || p(\theta)) = - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{All bivariate normals}\}$

What is  $q^*(\theta) = \operatorname{argmin}_{q \in \mathcal{Q}} \text{KL}(q(\theta) || p(\theta))$ ?



**Sufficiently expressive families recover the target distribution.**

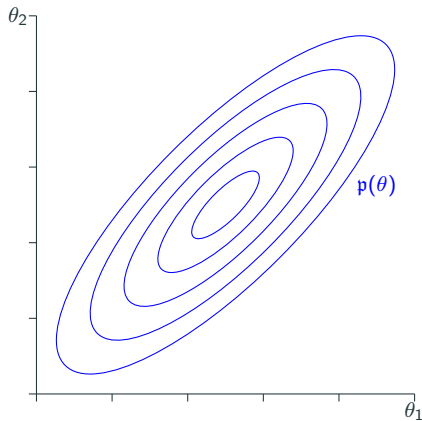
## KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



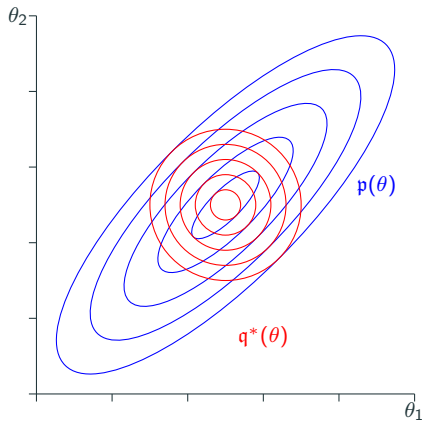
## KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



**KL minimizers “fit inside” the second argument.**

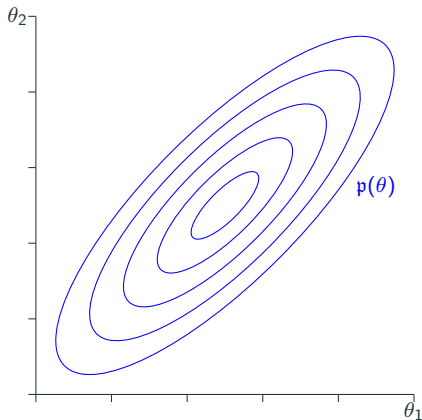
## KL divergence exercises

$$\begin{aligned} \text{KL}(p(\theta) || q(\theta)) = \\ - \mathbb{E}_{p(\theta)} [\log q(\theta)] + \mathbb{E}_{p(\theta)} [\log p(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(p(\theta) || q(\theta))$ ?



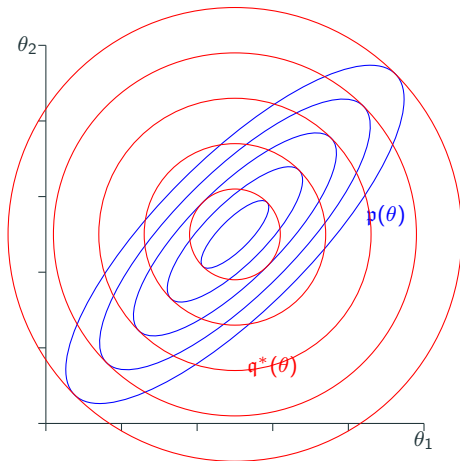
## KL divergence exercises

$$\begin{aligned} \text{KL}(p(\theta) || q(\theta)) = \\ - \mathbb{E}_{p(\theta)} [\log q(\theta)] + \mathbb{E}_{p(\theta)} [\log p(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(p(\theta) || q(\theta))$ ?



**KL minimizers “fit inside” the second argument.**



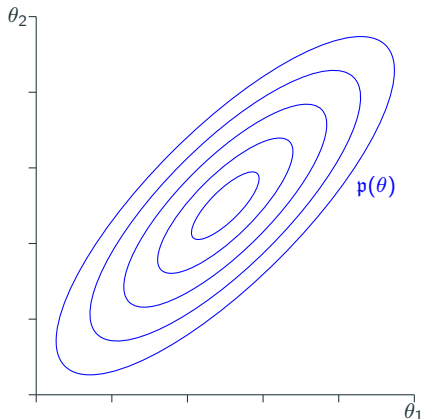
# KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$$\mathcal{Q} = \{\text{Bivariate normals}\}$$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \left( - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$ ?



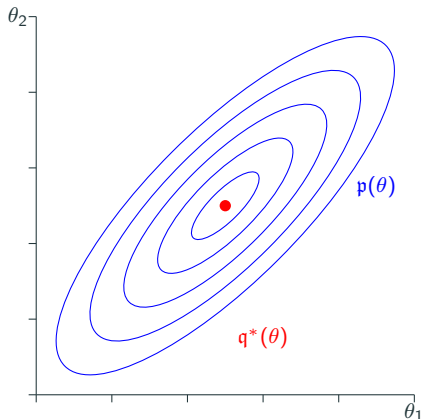
## KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Bivariate normals}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \left( - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$ ?



**Without the entropy, the KL minimizer concentrates on the maximum of  $\log p(\theta)$ .**

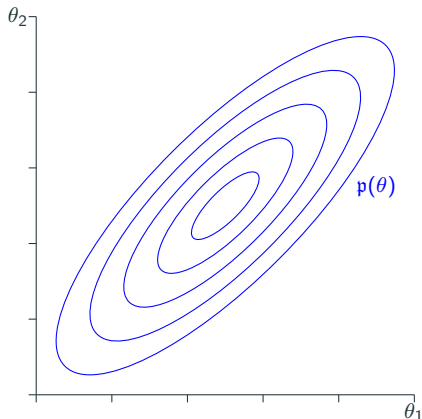
# KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$$\mathcal{Q} = \{\text{Bivariate normals}\}$$

What is  $q^*(\theta) =$   
 $\underset{q \in \mathcal{Q}}{\text{argmin}} \left( - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right) ?$



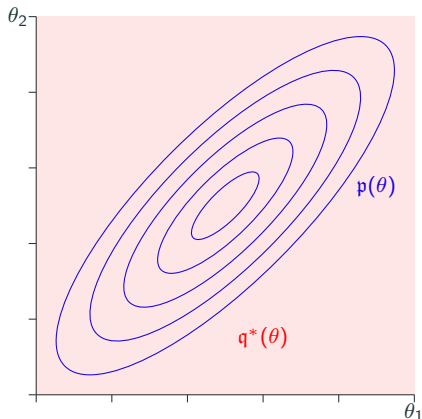
# KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Bivariate normals}\}$

What is  $q^*(\theta) =$   
 $\underset{q \in \mathcal{Q}}{\text{argmin}} \left( - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$ ?



**Without  $\log p(\theta)$ , the KL minimizer is infinitely dispersed.**

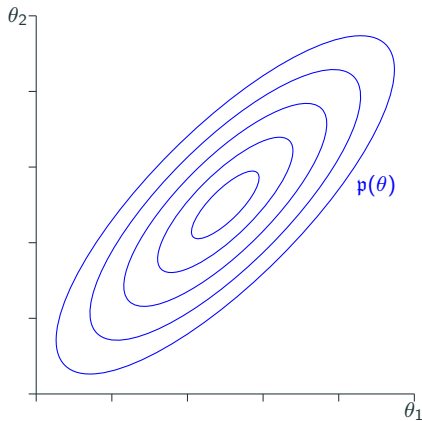
# KL divergence exercises

$$\text{KL}(q(\theta) || p(\theta)) = -\mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Point masses}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



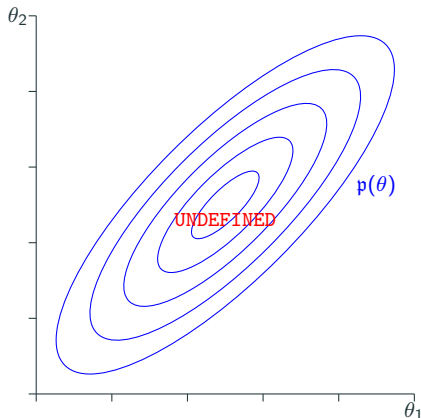
## KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{Point masses}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



**Without a common dominating measure, the KL divergence is undefined.**

## KL divergence exercises

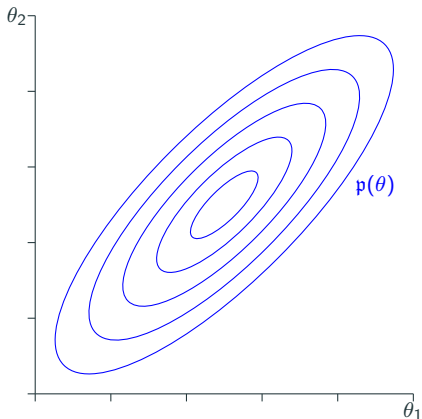
$$\text{KL}(q(\theta) || p(\theta)) =$$

$$- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{BVN with small, fixed variance}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



## KL divergence exercises

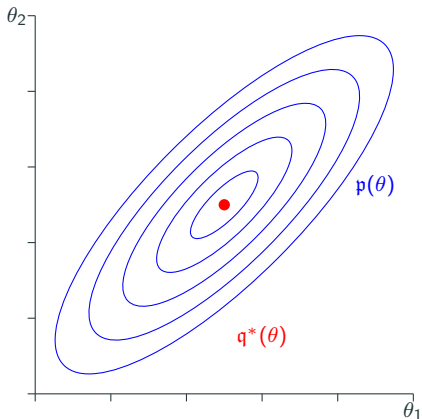
$$\text{KL}(q(\theta) || p(\theta)) =$$

$$- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$  = Correlated bivariate normal

$\mathcal{Q} = \{\text{BVN with small, fixed variance}\}$

What is  $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q(\theta) || p(\theta))$ ?



**Sufficiently concentrated distributions with constant entropy act like a point mass at the maximum of  $\log p(\theta)$ .**



## Recall the EM algorithm

Observations:  $y = (y_1, \dots, y_N)$

Unknown latent variables:  $z = (z_1, \dots, z_N)$

Unknown global parameter:  $\theta \in \mathbb{R}^D$ . We want:  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(y|\theta)$ .

## Recall the EM algorithm

Observations:  $y = (y_1, \dots, y_N)$

Unknown latent variables:  $z = (z_1, \dots, z_N)$

Unknown global parameter:  $\theta \in \mathbb{R}^D$ . We want:  $\hat{\theta} = \operatorname{argmax}_{\theta} \log p(y|\theta)$ .

---

The EM algorithm alternates between two steps. Starting at an iterate  $\hat{\theta}_{(i)}$ , repeat until convergence:

**The E-step:** Compute  $Q_{(i)}(\theta) := \mathbb{E}_{p(z|y, \hat{\theta}_{(i)})} [\log p(y|\theta, z) + \log p(z|\theta)]$

**The M-step:** Compute the next iterate  $\hat{\theta}_{(i+1)} := \operatorname{argmax}_{\theta} Q_{(i)}(\theta)$

## Recall the EM algorithm

Observations:  $y = (y_1, \dots, y_N)$

Unknown latent variables:  $z = (z_1, \dots, z_N)$

Unknown global parameter:  $\theta \in \mathbb{R}^D$ . We want:  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(y|\theta)$ .

---

The EM algorithm alternates between two steps. Starting at an iterate  $\hat{\theta}_{(i)}$ , repeat until convergence:

**The E-step:** Compute  $Q_{(i)}(\theta) := \underset{p(z|y, \hat{\theta}_{(i)})}{\mathbb{E}} [\log p(y|\theta, z) + \log p(z|\theta)]$

**The M-step:** Compute the next iterate  $\hat{\theta}_{(i+1)} := \underset{\theta}{\operatorname{argmax}} Q_{(i)}(\theta)$

---

The EM algorithm works / is useful when:

- The joint log probability  $\log p(y|\theta, z) + \log p(z|\theta)$  is easy to write down
- The posterior  $p(z|y, \theta)$  is easy to compute
- The marginalizing integral  $p(y|\theta) = \int p(y|\theta, z)p(z|\theta)dz$  is hard

## Recall the EM algorithm

Observations:  $y = (y_1, \dots, y_N)$

Unknown latent variables:  $z = (z_1, \dots, z_N)$

Unknown global parameter:  $\theta \in \mathbb{R}^D$ . We want:  $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(y|\theta)$ .

---

The EM algorithm alternates between two steps. Starting at an iterate  $\hat{\theta}_{(i)}$ , repeat until convergence:

**The E-step:** Compute  $Q_{(i)}(\theta) := \underset{p(z|y, \hat{\theta}_{(i)})}{\mathbb{E}} [\log p(y|\theta, z) + \log p(z|\theta)]$

**The M-step:** Compute the next iterate  $\hat{\theta}_{(i+1)} := \underset{\theta}{\operatorname{argmax}} Q_{(i)}(\theta)$

---

The EM algorithm works / is useful when:

- The joint log probability  $\log p(y|\theta, z) + \log p(z|\theta)$  is easy to write down
  - The posterior  $p(z|y, \theta)$  is easy to compute
  - The marginalizing integral  $p(y|\theta) = \int p(y|\theta, z)p(z|\theta)dz$  is hard
- 

Is the EM algorithm VI?

Can you spot the lie?

## The EM algorithm as VI

Let  $\mathcal{Q}_z$  denote a family of distributions on  $z$ , parameterized by a finite-dimensional parameter  $\eta$ , such that  $p(z|\theta, y) \in \mathcal{Q}_z$  for the observed  $y$  and all  $\theta$ .

**Exercise:** When does  $\mathcal{Q}_z$  exist? (Indexed by a finite-dimensional parameter  $\eta$ .)

## The EM algorithm as VI

Let  $\mathcal{Q}_z$  denote a family of distributions on  $z$ , parameterized by a finite-dimensional parameter  $\eta$ , such that  $p(z|\theta, y) \in \mathcal{Q}_z$  for the observed  $y$  and all  $\theta$ .

**Exercise:** When does  $\mathcal{Q}_z$  exist? (Indexed by a finite-dimensional parameter  $\eta$ .)

Let  $q(z|\hat{\eta}(\theta)) := p(z|\theta, y)$ .

In an abuse of notation, write  $\eta \in \mathcal{Q}_z$  for  $\eta \in \{\eta : q(z|\eta) \in \mathcal{Q}_z\}$ .

## The EM algorithm as VI

Let  $\mathcal{Q}_z$  denote a family of distributions on  $z$ , parameterized by a finite-dimensional parameter  $\eta$ , such that  $p(z|\theta, y) \in \mathcal{Q}_z$  for the observed  $y$  and all  $\theta$ .

**Exercise:** When does  $\mathcal{Q}_z$  exist? (Indexed by a finite-dimensional parameter  $\eta$ .)

Let  $q(z|\hat{\eta}(\theta)) := p(z|\theta, y)$ .

In an abuse of notation, write  $\eta \in \mathcal{Q}_z$  for  $\eta \in \{\eta : q(z|\eta) \in \mathcal{Q}_z\}$ .

Then:

$$\begin{aligned}\log p(y|\theta) &= \log p(y|\theta) + \text{KL}(q(z|\hat{\eta}(\theta))||p(z|\theta, y)) \\ &= \log p(y|\theta) + \underset{\eta \in \mathcal{Q}_z}{\operatorname{argmax}} (-\text{KL}(q(z|\eta)||p(z|\theta, y))) \quad \star \\ &= \underset{\eta \in \mathcal{Q}_z}{\operatorname{argmax}} (\log p(y|\theta) - \text{KL}(q(z|\eta)||p(z|\theta, y))) \\ &= \underset{\eta \in \mathcal{Q}_z}{\operatorname{argmax}} \left( \log p(y|\theta) + \mathbb{E}_{q(z|\eta)} [\log p(z|\theta, y) - \log q(z|\eta)] \right) \\ &= \underset{\eta \in \mathcal{Q}_z}{\operatorname{argmax}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y|\theta) + \log p(z|\theta, y) - \log q(z|\eta)] \right) \\ &= \underset{\eta \in \mathcal{Q}_z}{\operatorname{argmax}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] - \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \quad \star\star\end{aligned}$$

# The EM algorithm as VI

From the previous slide, the marginal MLE is given by

$$\hat{\theta} := \operatorname{argmax}_{\theta} \log p(y|\theta)$$

$$= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \quad \star$$

$$= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \quad \star\star$$



# The EM algorithm as VI

From the previous slide, the marginal MLE is given by

$$\begin{aligned}\hat{\theta} &:= \operatorname{argmax}_{\theta} \log p(y|\theta) \\ &= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \quad \star \\ &= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \quad \star\star\end{aligned}$$

---

**The EM algorithm revisited.** Starting at an iterate  $\hat{\theta}_{(i)}$ :

**The E-step:**

1. For a fixed  $\hat{\theta}_{(i)}$ , optimize  $\star$  for  $\eta$ . Since only the KL divergence depends on  $\eta$ , the optimum is  $\hat{\eta}(\hat{\theta}_{(i)})$ , and  $q(z|\hat{\eta}(\hat{\theta}_{(i)})) = p(z|\hat{\theta}_{(i)}, y)$ .
2. Then use  $\hat{\eta}(\hat{\theta}_{(i)})$  to compute the expectation in  $\star\star$  as a function of  $\theta$ .

**The M-step:** Keeping  $\eta$  fixed at  $\hat{\eta}(\hat{\theta}_{(i)})$ , optimize  $\star\star$  as a function of  $\theta$  to give  $\hat{\theta}_{i+1}$ .

The entropy  $\mathbb{E}_{q(z|\eta)} [\log q(z|\eta)]$  does not depend on  $\theta$  and can be ignored.

# The EM algorithm as VI

From the previous slide, the marginal MLE is given by

$$\begin{aligned}\hat{\theta} &:= \operatorname{argmax}_{\theta} \log p(y|\theta) \\ &= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \quad \star \\ &= \operatorname{argmax}_{\theta} \operatorname{argmax}_{\eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \quad \star\star\end{aligned}$$

---

**The EM algorithm revisited.** Starting at an iterate  $\hat{\theta}_{(i)}$ :

**The E-step:**

1. For a fixed  $\hat{\theta}_{(i)}$ , optimize  $\star$  for  $\eta$ . Since only the KL divergence depends on  $\eta$ , the optimum is  $\hat{\eta}(\hat{\theta}_{(i)})$ , and  $q(z|\hat{\eta}(\hat{\theta}_{(i)})) = p(z|\hat{\theta}_{(i)}, y)$ .
2. Then use  $\hat{\eta}(\hat{\theta}_{(i)})$  to compute the expectation in  $\star\star$  as a function of  $\theta$ .

**The M-step:** Keeping  $\eta$  fixed at  $\hat{\eta}(\hat{\theta}_{(i)})$ , optimize  $\star\star$  as a function of  $\theta$  to give  $\hat{\theta}_{i+1}$ . The entropy  $\mathbb{E}_{q(z|\eta)} [\log q(z|\eta)]$  does not depend on  $\theta$  and can be ignored.

---

$\Rightarrow$  The EM algorithm is coordinate ascent on the objective

$$f(\theta, \eta) = \log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y)) .$$

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

**Corollaries:**

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

**Corollaries:**

- The EM algorithm converges to a local optimum of  $\log p(y|\theta)$ .

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

**Corollaries:**

- The EM algorithm converges to a local optimum of  $\log p(y|\theta)$ .
- The EM algorithm is VI, and you don't need to optimize with coordinate ascent.

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

**Corollaries:**

- The EM algorithm converges to a local optimum of  $\log p(y|\theta)$ .
- The EM algorithm is VI, and you don't need to optimize with coordinate ascent.
- If both  $p(z|\theta, y)$  and  $p(z, y|\theta)$  are easy, then so is  $p(y|\theta)$ . (This was the lie.)

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

## Corollaries:

- The EM algorithm converges to a local optimum of  $\log p(y|\theta)$ .
- The EM algorithm is VI, and you don't need to optimize with coordinate ascent.
- If both  $p(z|\theta, y)$  and  $p(z, y|\theta)$  are easy, then so is  $p(y|\theta)$ . (This was the lie.)
  - **Exercise:** Prove this a different way using exponential families (in real-life problems,  $p(z|y, \theta)$  is only really tractable in exponential families).

$$\hat{\theta}, \hat{\eta} := \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))).$$

The EM algorithm is coordinate ascent on the preceding objective.  
[Neal and Hinton, 1998]

## Corollaries:

- The EM algorithm converges to a local optimum of  $\log p(y|\theta)$ .
- The EM algorithm is VI, and you don't need to optimize with coordinate ascent.
- If both  $p(z|\theta, y)$  and  $p(z, y|\theta)$  are easy, then so is  $p(y|\theta)$ . (This was the lie.)
  - **Exercise:** Prove this a different way using exponential families (in real-life problems,  $p(z|y, \theta)$  is only really tractable in exponential families).
- If  $p(z|\theta, y)$  is intractable, we can now consider different approximating families which may not contain  $p(z|\theta, y)$ .



## Different approximating families: Point masses on $z$ .

Suppose instead of  $Q_z$  we used  $Q_z^{pm}$ , a family of constant-entropy near-point mass distributions on  $z$ , located at some free parameter  $\eta$ .

## Different approximating families: Point masses on $z$ .

Suppose instead of  $\mathcal{Q}_z$  we used  $\mathcal{Q}_z^{pm}$ , a family of constant-entropy near-point mass distributions on  $z$ , located at some free parameter  $\eta$ .

Then

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] \right) \\ &= \operatorname{argmax}_{\theta, z} \log p(y, z|\theta). \end{aligned}$$

## Different approximating families: Point masses on $z$ .

Suppose instead of  $\mathcal{Q}_z$  we used  $\mathcal{Q}_z^{pm}$ , a family of constant-entropy near-point mass distributions on  $z$ , located at some free parameter  $\eta$ .

Then

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] \right) \\ &= \operatorname{argmax}_{\theta, z} \log p(y, z|\theta). \end{aligned}$$

$\Rightarrow$  The Neyman-Scott paradox occurs because point masses are poor approximations for the distribution  $p(z|\theta, y)$ .

## Different approximating families: Point masses on $z$ .

Suppose instead of  $\mathcal{Q}_z$  we used  $\mathcal{Q}_z^{pm}$ , a family of constant-entropy near-point mass distributions on  $z$ , located at some free parameter  $\eta$ .

Then

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] + \mathbb{E}_{q(z|\eta)} [\log q(z|\eta)] \right) \\ &= \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z^{pm}} \left( \mathbb{E}_{q(z|\eta)} [\log p(y, z|\theta)] \right) \\ &= \operatorname{argmax}_{\theta, z} \log p(y, z|\theta). \end{aligned}$$

**$\Rightarrow$  The Neyman-Scott paradox occurs because point masses are poor approximations for the distribution  $p(z|\theta, y)$ .**

**Exercise:** Recall that the Neyman-Scott paradox disappears when, instead of pairs, we have many observations, all from the same  $z_n$ . Can you use the VI perspective on the marginal EM algorithm to explain this phenomenon?

## Different approximating families: Point masses on $\theta$ .

Let  $\mathcal{Q}_{\theta}^{pm}$  denote a family of constant-entropy near-point mass distributions on  $\theta$ , located at some free parameter  $\vartheta$ . Assume a uniform prior on  $\theta$ .

## Different approximating families: Point masses on $\theta$ .

Let  $\mathcal{Q}_\theta^{pm}$  denote a family of constant-entropy near-point mass distributions on  $\theta$ , located at some free parameter  $\vartheta$ . Assume a uniform prior on  $\theta$ .

Then:

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(\theta|\vartheta)} [\log p(y|\theta)] - \mathbb{E}_{q(\theta|\vartheta)} [\log q(\theta|\vartheta)] - \text{KL}(q(z|\eta) || p(z|\theta, y)) \right) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \text{KL}(q(\theta|\vartheta)q(z|\eta) || \log p(z, \theta|y)). \end{aligned}$$

## Different approximating families: Point masses on $\theta$ .

Let  $\mathcal{Q}_\theta^{pm}$  denote a family of constant-entropy near-point mass distributions on  $\theta$ , located at some free parameter  $\vartheta$ . Assume a uniform prior on  $\theta$ .

Then:

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(\theta|\vartheta)} [\log p(y|\theta)] - \mathbb{E}_{q(\theta|\vartheta)} [\log q(\theta|\vartheta)] - \text{KL}(q(z|\eta) || p(z|\theta, y)) \right) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \text{KL}(q(\theta|\vartheta)q(z|\eta) || \log p(z, \theta|y)). \end{aligned}$$

$\Rightarrow$  **The marginal MLE is a point-mass approximation to the posterior with a uniform prior.** It will be a good approximation when  $p(\theta|y)$  is approximately a point mass.

## Different approximating families: Point masses on $\theta$ .

Let  $\mathcal{Q}_\theta^{pm}$  denote a family of constant-entropy near-point mass distributions on  $\theta$ , located at some free parameter  $\vartheta$ . Assume a uniform prior on  $\theta$ .

Then:

$$\begin{aligned} & \operatorname{argmax}_{\theta, \eta \in \mathcal{Q}_z} (\log p(y|\theta) - \text{KL}(q(z|\eta) || p(z|\theta, y))) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \left( \mathbb{E}_{q(\theta|\vartheta)} [\log p(y|\theta)] - \mathbb{E}_{q(\theta|\vartheta)} [\log q(\theta|\vartheta)] - \text{KL}(q(z|\eta) || p(z|\theta, y)) \right) \\ &= \operatorname{argmax}_{\vartheta \in \mathcal{Q}_\theta^{pm}, \eta \in \mathcal{Q}_z} \text{KL}(q(\theta|\vartheta)q(z|\eta) || \log p(z, \theta|y)). \end{aligned}$$

$\Rightarrow$  **The marginal MLE is a point-mass approximation to the posterior with a uniform prior.** It will be a good approximation when  $p(\theta|y)$  is approximately a point mass.

**Exercise:** Under what circumstances is  $p(\theta|y)$  is approximately a point mass?



## Different approximating families.

Suppose we can't compute  $p(z|\theta, y)$  and / or we think that  $p(\theta|y)$  may not be well-approximated by a point mass.

Choose some tractable approximating family  $q(\theta, z|\gamma) \in \mathcal{Q}_{\theta z}$ . Then find

$$\hat{\gamma} := \operatorname{argmin}_{\gamma \in \mathcal{Q}_{\theta z}} \operatorname{KL}(q(\theta, z|\gamma) || p(\theta, z|y)).$$

## Different approximating families.

Suppose we can't compute  $p(z|\theta, y)$  and / or we think that  $p(\theta|y)$  may not be well-approximated by a point mass.

Choose some tractable approximating family  $q(\theta, z|\gamma) \in \mathcal{Q}_{\theta z}$ . Then find

$$\hat{\gamma} := \operatorname{argmin}_{\gamma \in \mathcal{Q}_{\theta z}} \operatorname{KL}(q(\theta, z|\gamma) || p(\theta, z|y)).$$

Now we're doing **“Variational Bayes” (VB)**.

## Different approximating families.

Suppose we can't compute  $p(z|\theta, y)$  and / or we think that  $p(\theta|y)$  may not be well-approximated by a point mass.

Choose some tractable approximating family  $q(\theta, z|\gamma) \in \mathcal{Q}_{\theta z}$ . Then find

$$\hat{\gamma} := \operatorname{argmin}_{\gamma \in \mathcal{Q}_{\theta z}} \operatorname{KL}(q(\theta, z|\gamma) || p(\theta, z|y)).$$

**Now we're doing “Variational Bayes” (VB).**

The EM algorithm — and, indeed, the MLE — can be understood as Variational Bayes with a uniform prior and particular choices of approximating distributions.

## Different approximating families.

Some common approximating families:

- Factorizing families, e.g.  $q(\theta, z|\gamma) = q(\theta|\gamma)q(z|\gamma)$ . These families model some components of the posterior as independent.
  - For historical reasons, this is known as a **mean-field approximation** [Wainwright and Jordan, 2008].
- Factorizing families + an exponential family assumption.
- Normal approximations (possibly after an invertible unconstraining transformation):  $q(\theta, z|\gamma) = \mathcal{N}(\theta, z|\gamma)$ .
- Independent normal approximations. This is used by a lot of “black-box VI” methods [Ranganath et al., 2014, Kucukelbir et al., 2017].

## Different approximating families.

Some common approximating families:

- Factorizing families, e.g.  $q(\theta, z|\gamma) = q(\theta|\gamma)q(z|\gamma)$ . These families model some components of the posterior as independent.
  - For historical reasons, this is known as a **mean-field approximation** [Wainwright and Jordan, 2008].
- Factorizing families + an exponential family assumption.
- Normal approximations (possibly after an invertible unconstraining transformation):  $q(\theta, z|\gamma) = \mathcal{N}(\theta, z|\gamma)$ .
- Independent normal approximations. This is used by a lot of “black-box VI” methods [Ranganath et al., 2014, Kucukelbir et al., 2017].

---

What do you need from an approximating family? Expressivity, plus:

$$\text{KL}(q||p) := \underbrace{\mathbb{E}_{q(\xi|\eta)} [\log q(\xi|\eta)]}_{\text{Tractable entropy}} - \underbrace{\mathbb{E}_{q(\xi|\eta)} [\log p(\xi)]}_{\text{Tractable expectations}}$$

# Different approximating families.

Some common approximating families:

- Factorizing families, e.g.  $q(\theta, z|\gamma) = q(\theta|\gamma)q(z|\gamma)$ . These families model some components of the posterior as independent.
    - For historical reasons, this is known as a **mean-field approximation** [Wainwright and Jordan, 2008].
  - Factorizing families + an exponential family assumption.
  - Normal approximations (possibly after an invertible unconstraining transformation):  $q(\theta, z|\gamma) = \mathcal{N}(\theta, z|\gamma)$ .
  - Independent normal approximations. This is used by a lot of “black-box VI” methods [Ranganath et al., 2014, Kucukelbir et al., 2017].
- 

What do you need from an approximating family? Expressivity, plus:

$$\text{KL}(q||p) := \underbrace{\mathbb{E}_{q(\xi|\eta)} [\log q(\xi|\eta)]}_{\text{Tractable entropy}} - \underbrace{\mathbb{E}_{q(\xi|\eta)} [\log p(\xi)]}_{\text{Tractable expectations}}$$

- Monte Carlo is often used for the expectations.
  - See, e.g., Ranganath et al. [2014].
- The entropy is harder. In general, there is a tradeoff between expressivity and tractable entropy.
  - “Normalizing flows” are an example of a highly expressive approximating family (neural nets!) designed to maintain a tractable entropy. [Rezende and Mohamed, 2015]

## VI in practice: The Criteo dataset

As an example application of VB, consider a logistic regression with random effects fit (generalized linear mixed model) to an internet advertising dataset from Criteo Labs with  $N = 61895$  datapoints [Giordano et al., 2018, Section 5.3].

We want to estimate:

$\beta$  : Regression parameters (5-dimensional)

$u$  : Random effects (5000-dimensional)

$\mu$  : Random effect mean (intercept)

$\tau$  : Random effect precision (inverse variance).

## VI in practice: The Criteo dataset

As an example application of VB, consider a logistic regression with random effects fit (generalized linear mixed model) to an internet advertising dataset from Criteo Labs with  $N = 61895$  datapoints [Giordano et al., 2018, Section 5.3].

We want to estimate:

- $\beta$  : Regression parameters (5-dimensional)
- $u$  : Random effects (5000-dimensional)
- $\mu$  : Random effect mean (intercept)
- $\tau$  : Random effect precision (inverse variance).

We use the following VB approximation:

$$\begin{aligned}q(\beta_k) &= \mathcal{N}(\beta_k; \eta_{\beta_k}), \text{ for } k = 1, \dots, K_x \\q(u_t) &= \mathcal{N}(u_t; \eta_{u_t}), \text{ for } t = 1, \dots, T \\q(\tau) &= \text{Gamma}(\tau; \eta_\tau) \\q(\mu) &= \mathcal{N}(\mu; \eta_\mu) \\q(\theta) &= q(\tau) q(\mu) \prod_{k=1}^{K_x} q(\beta_k) \prod_{t=1}^T q(u_t).\end{aligned}$$

We will compare the joint MAP ( $\approx$  MLE), MCMC, and the VB approximation.



# VI in practice: The Criteo dataset

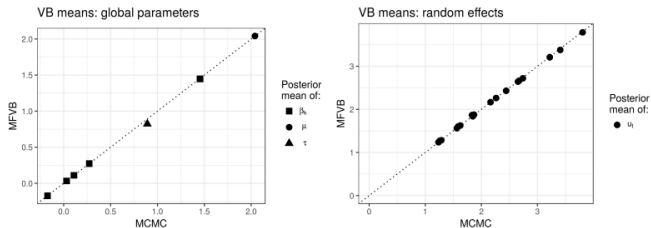


Figure 13: Comparison of MCMC and MFVB means

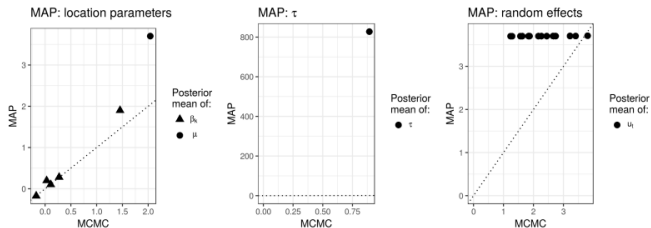
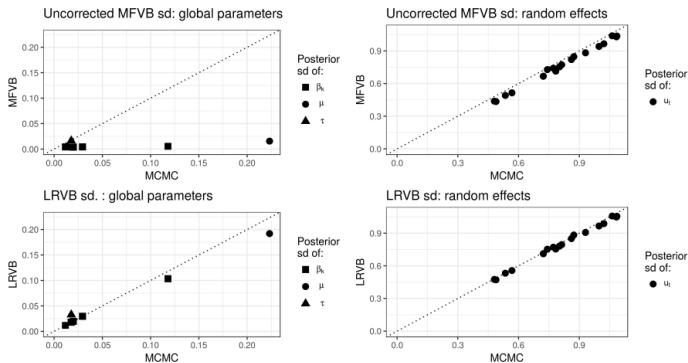


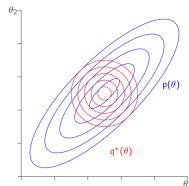
Figure 14: Comparison of MCMC and Laplace means

# VI in practice: The Criteo dataset



Note that standard mean-field VB under-estimates posterior covariances. We have a paper about alleviating this problem using “linear response.” [Giordano et al., 2018]

The Hessian was singular at the MAP, so the Laplace approximation could not be computed.



VB is slower than the MAP, but much faster than MCMC.

Method	Seconds
MAP (optimum only)	12
VB (optimum only)	57
VB (including sensitivity for $\beta$ )	104
VB (including sensitivity for $\beta$ and $u$ )	553
MCMC (Stan)	21066

## VI in practice: Additional results

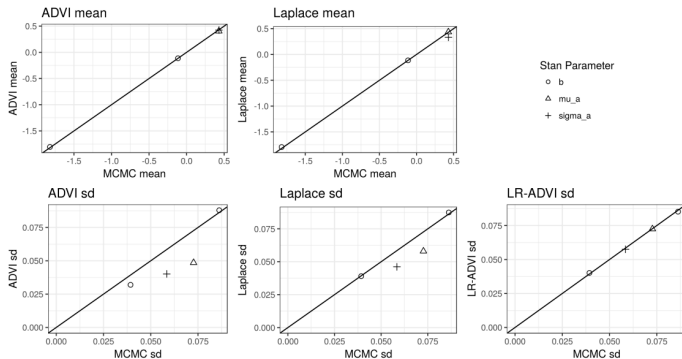


Figure 7: Election model

## VI in practice: Additional results

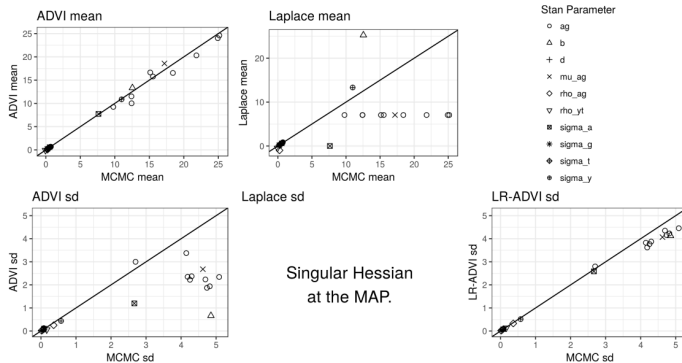


Figure 8: Sesame Street model

## VI in practice: Additional results

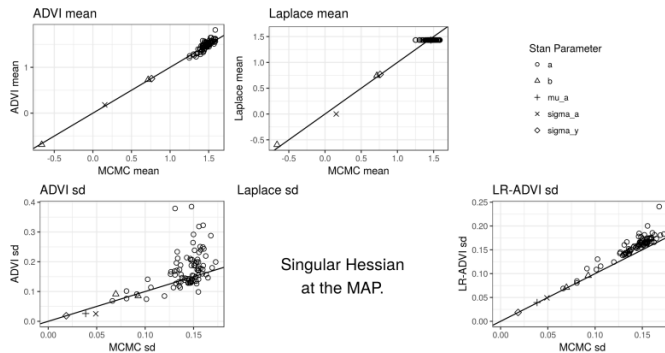


Figure 9: Radon model

## VI in practice: Additional results

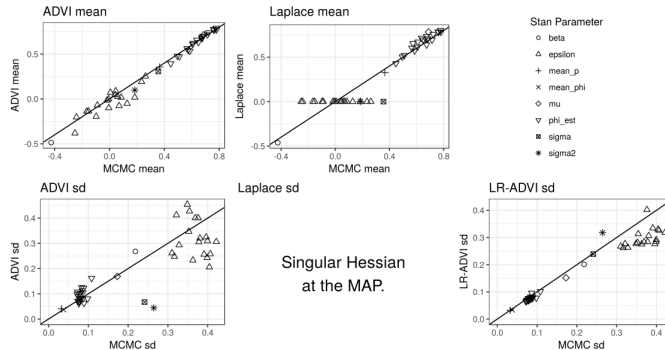


Figure 10: Ecology model

## VI in practice: Additional results

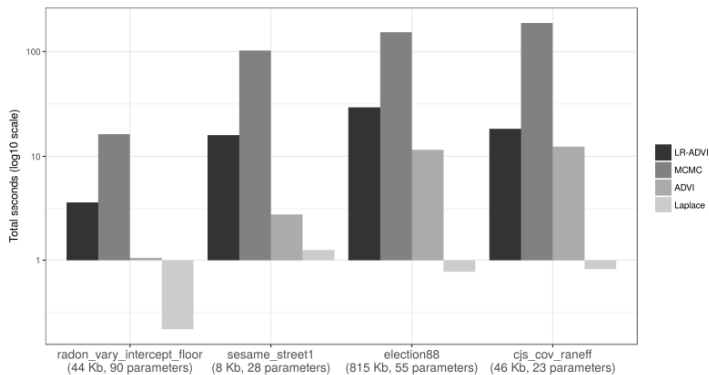


Figure 11: Comparison of timing in ADVI experiments



- Variational inference can be thought of as *approximate marginalization using optimization*.
- Not only is variational inference closely related to familiar existing frequentist methods (like the EM algorithm), it can help us understand those methods better.
- The key to a good variational approximation is *tractable expectations*, *tractable entropy*, and *expressivity*.
- It is important to be aware of variational inference's shortcomings (e.g., underestimation of variance in the mean field approximation).
- There are a zillion topics to work on in variational inference. A good place to start reading is Blei et al. [2017].

**Thanks for having me!**

- Luca Ambrogioni, Umut Güçlü, Yağmur Güçlütürk, Max Hinde, Eric Maris, and Marcel AJ van Gerven. Wasserstein variational inference. *arXiv preprint arXiv:1805.11284*, 2018.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Alison L Gibbs and Francis Edward Su. On choosing and bounding probability metrics. *International statistical review*, 70(3):419–435, 2002.
- Ryan Giordano, Tamara Broderick, and Michael I Jordan. Covariances, robustness and variational bayes. *Journal of Machine Learning Research*, 19(51), 2018.
- Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M Blei. Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Yingzhen Li and Richard E Turner. Variational inference with rényi divergence. *stat*, 1050:6, 2016.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. *arXiv preprint arXiv:1608.04471*, 2016.
- Radford M Neal and Geoffrey E Hinton. A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer, 1998.
- Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial intelligence and statistics*, pages 814–822. PMLR, 2014.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- Martin J Wainwright and Michael I. Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.